

# Forecasting the 2016-2017 Central Apennines Earthquake Sequence with a Neural Point Process

Samuel Stockman <sup>1</sup>, Daniel J. Lawson <sup>1</sup>, Maximilian J. Werner <sup>2</sup>

<sup>1</sup>School of Mathematics, University of Bristol

<sup>2</sup>School of Earth Sciences, University of Bristol

## Key Points:

- We construct a new machine learning variant of point processes for short-term earthquake forecasting enhanced catalogs.
- The neural point process gains higher forecasting performance from the low magnitude data than ETAS and is faster to train.
- This forecasting performance on the 2016 Central Italy sequence motivates continued development in this class of models.

## Abstract

Point processes have been dominant in modeling the evolution of seismicity for decades, with the Epidemic-Type Aftershock Sequence (ETAS) model being most popular. Recent advances in machine learning have constructed highly flexible point process models using neural networks to improve upon existing parametric models. We investigate whether these flexible point process models can be applied to short-term seismicity forecasting by extending an existing temporal neural model to the magnitude domain and we show how this model can forecast earthquakes above a target magnitude threshold. We first demonstrate that the neural model can fit synthetic ETAS data, however, requiring less computational time because it is not dependent on the full history of the sequence. By artificially emulating short-term aftershock incompleteness in the synthetic dataset, we find that the neural model outperforms ETAS. Using a new enhanced catalog from the 2016-2017 Central Apennines earthquake sequence, we investigate the predictive skill of ETAS and the neural model with respect to the lowest input magnitude. Constructing multiple forecasting experiments using the Visso, Norcia and Campotosto earthquakes to partition training and testing data, we target M3+ events. We find both models perform similarly at previously explored thresholds (e.g., above M3), but lowering the threshold to M1.2 reduces the performance of ETAS unlike the neural model. We argue that some of these gains are due to the neural model’s ability to handle incomplete data. The robustness to missing data and speed to train the neural model present it as an encouraging competitor in earthquake forecasting.

## Plain Language Summary

For decades, the Epidemic-Type Aftershock Sequence (ETAS) model has been the most popular way of forecasting earthquakes over short time spans (days/weeks). It is formulated mathematically as a point process, a general class of statistical model describing the random occurrence of points in time. Recently the machine learning community have used neural networks to make point processes more expressive and titled them neural point processes. In this study we investigate whether a neural point process can compete with the ETAS model. We find that the two models perform similarly on computer simulated data; however, the neural model is much faster with large datasets and is not hindered if there is missing data for smaller earthquakes. Most earthquake catalogs contain missing data due to varying capability in our detection methods, therefore we need

models that are robust to this missingness. We then find that the neural model outperforms ETAS on a new catalog for the 2016-2017 Central Apennines earthquake sequence, which through machine learning detection contains thousands of previously undetected small magnitude events. We argue that some of this improvement can in fact be explained by missing data. These results present neural point processes as an encouraging competitor in earthquake forecasting.

## 1 Introduction

The construction of machine learning algorithms for detecting the arrival times of earthquake phases (eg. Zhu and Beroza (2019)) combined with an accelerated growth in the number of seismic sensors, has meant that sizes of earthquake catalogs have grown substantially (Kong et al., 2019). With the amount of available seismicity data increasing, current forecasting methods that include the full history of the catalog in the form of all event pairs are increasingly inefficient and might not be flexible enough to incorporate this additional data, thus the need for the application of methods developed in the machine learning community is becoming more apparent. However, there exists a disconnect between the tools used by statistical seismologists and those in the machine learning community that apply their methods to seismic data. This work attempts to bridge that gap (to some extent) by considering a machine learning variant of point processes. Point processes are a class of models that contain the Epidemic-Type Aftershock Sequence (ETAS) model, a widely accepted and used point process model for earthquakes (Ogata, 1988, 1998; Marzocchi et al., 2014; Mancini et al., 2019). In working with a machine learning method with a conditional intensity function (the function that explicitly defines a point process) (Zhuang et al., 2012), we present a model that is directly comparable to ETAS models, the current benchmark for short-term earthquake forecasting, yet now with desirable properties such as flexibility and scalability. The machine learning variant of point processes we introduce are known as neural point processes.

At the heart of neural point processes is the use of a recurrent neural network (RNN) to learn a compact representation of the history of events, first introduced by Du et al. (2016). The sequential nature of the way data pass through RNNs makes them an ideal modeling tool for temporal data. Instead of directly summing over all past events, as in models based on the Hawkes process (Hawkes, 1971), a fixed length vector representation of the past is learnt and updated at each new time step. Forecasts can then be made

by modelling the conditional intensity function on this vector (Xiao et al., 2017; Li et al., 2018; Upadhyay et al., 2018; Huang et al., 2019; Omi et al., 2019), or instead through directly modelling the probability of the next event (Shchur et al., 2019). We direct the reader to Shchur et al. (2021) for a review of neural point process models. None of these models, however, are directly suitable for describing the temporal behaviour of seismicity, which includes a continuous mark space for the magnitude of each earthquake as well as the times. We require a model that is dependent on continuous marks as well as forecasts them. To the best of our knowledge, neither of these requirements are satisfied in the existing temporal point process literature.

In this work we extend the architecture introduced by Omi et al. (2019) so that it may also deal with earthquake magnitudes. For this we require a model that is dependent on previous event magnitudes, can forecast subsequent magnitudes as well as forecast earthquakes above a threshold magnitude. Distinguishing between a threshold for the input magnitude and a threshold for the target earthquakes is a problem specific requirement for earthquake forecasting, so does not exist in other works on temporal point processes. We choose to extend Omi et al. (2019) to give the most flexible representation of the intensity, since they use a fully non-parametric approach compared to other intensity based methods that use a semi-parametric approach. Working with the intensity function rather than directly modelling the likelihood of the next event provides a model that is closer in interpretation to ETAS and provides a natural way to forecast earthquakes above some target threshold magnitude, detailed in section 4.2.

Although the use of the ETAS model (Ogata, 1988) has been the dominant way for modelling seismicity in both retrospective and fully prospective forecasting experiments (e.g. Woessner et al., 2011; Taroni et al., 2018; Cattania et al., 2018; Mancini et al., 2019, 2020) as well as in operational earthquake forecasting (Marzocchi et al., 2014; Rhoades et al., 2016; Field et al., 2017), it is restricted to its rigid parametric form. As ETAS only describes the self-exciting nature of seismicity, it cannot capture any kind of inhibition or release of stress such as captured by stress-release models (Zheng & Vere-Jones, 1991; Xiaogu & Vere-Jones, 1994; Bebbington & Harte, 2003) or models based on elastostatic stress transfer and Coulomb Rate-and-State (CRS) friction (Dieterich, 1994). Furthermore, foreshock activity that differs from ETAS has also been observed (McGuire et al., 2005; Brodsky, 2011; Lippiello et al., 2012; Ogata & Katsura, 2014). Beyond this understanding that ETAS is misspecified, there are also difficulties and inef-

iciencies with fitting and forecasting. To estimate the intensity, ETAS sums over all previous earthquakes, which requires substantial memory and slows the fitting process and forecasting simulations. For large earthquakes in the past this is important, because their contribution can last more than 100 years (Utsu et al., 1995). However, for smaller earthquakes particularly found in enhanced catalogs, one expects the contribution to be close to zero after a far shorter amount of time, making summing over these terms inefficient (Helmstetter & Sornette, 2003; Marsan & Lengline, 2008). Furthermore, a particular difficulty with fitting the ETAS model is that there needs to be a reliable estimate of the completeness across the time of the catalog and this needs to be incorporated into the model. Failing to do so will result in biases (Hainzl, 2016b; Zhuang et al., 2017; Seif et al., 2017). Methods that attempt to do this either truncate the periods of time where the catalog is most incomplete (Kagan, 1991; Hainzl et al., 2008), leading to parameters that can be dominated by a few aftershock sequences, or attempt to model the data incompleteness itself (Omi et al., 2014; Hainzl, 2016a, 2016b; Mizrahi et al., 2021b), but this adds additional computational requirements over a standard ETAS model.

To benchmark our proposed neural point process with ETAS, we design forecasting experiments on both synthetic data as well as a new enhanced catalog for the 2016–2017 Amatrice–Visso–Norcia (AVN) seismic sequence. The catalog generated by Tan et al. (2021b), containing roughly 900,000 earthquakes, was generated using a deep neural network based phase picker for earthquake arrival times (Zhu & Beroza, 2019). As a result, the size of the catalog increased 10 fold from the routine catalog generated by the Italian National Institute of Geophysics and Volcanology (INGV). The INGV catalog has been used in several retrospective forecasting experiments (Ebrahimian & Jalayer, 2017; Marzocchi et al., 2017; Mancini et al., 2019, 2022), but there has yet to be much development of forecasting models using enhanced catalogs such as this one and, given that they contain considerably more earthquakes, investigations into how we can harness these machine learning generated enhanced catalogs are essential. The AVN sequence contains ten  $M5+$  events during a five month period over an 80 km long normal-fault system (Mancini et al., 2019). The number of large earthquakes as well as the compactness of the region on which they occur make this sequence preferable for testing purely temporal forecasting models which contain no spatial covariates. We do still expect some loss of information by ignoring the spatial covariates, particularly for smaller earthquakes where there is a spatial extent across which earthquakes won’t interact.

We seek to understand how taking different magnitude thresholds and temporal partitions of our datasets affects the performance of the two models. Through altering these two aspects, we naturally change the amount of data shown to each model so that we may see how sensitive their forecasting performance is to training sample size (Wang et al., 2010). Through partitioning in time we can see how the performance is affected by the number of major earthquakes that each model is trained on. By altering the magnitude threshold of the input catalog, we seek to improve the predictive skill of forecasts by using the hypothesis that small earthquakes should help to forecast the moderate-to-large earthquakes. Either from a time-independent perspective where large earthquakes are found to nucleate in areas that have a large density of small events (Kafka & Walcott, 1998; Kafka & Levin, 2000), or in time-dependent forecasting (eg. ETAS and CRS) where earthquake triggering is believed to exist at all scales (Helmstetter & Sornette, 2003; Marsan, 2005; Nandan et al., 2016), reducing the threshold of the input catalog generally leads to improved forecasting performance of moderate-to-large events (Helmstetter et al., 2007; Werner et al., 2011; Helmstetter & Werner, 2014; Mancini et al., 2022).

Particularly, Mancini et al. (2022) consider the same sequence as this study, and compare forecasting results from models trained on several different enhanced catalogs (including the Tan et al. (2021b) catalog used in this study). They find that the forecasting of M3+ events by CRS and ETAS models is not improved by training on the enhanced catalogs. When using the same catalog as this study, they see the models increase in performance as the input magnitude threshold is lowered from M5 to M3, but at the lowest two thresholds M1 and M2, the performance of ETAS is worse than for M3+. Direct comparison of results, however, isn't possible as they report information gains that are for spatio-temporal forecasts using the Poisson assumption of earthquake rates in gridded space-time windows. This study hopes to provide some further insight into the performance of forecasting models using the low magnitude earthquakes found in this catalog and presents neural point processes as a competitive model using such events.

## 2 Data

To benchmark our neural point process against ETAS we conduct forecasting experiments on both real data from the Amatrice-Visso-Norcia sequence as well as synthetic data generated by ETAS. Since we are making comparisons about temporal models, in both catalogs, we remove all spatial covariates.

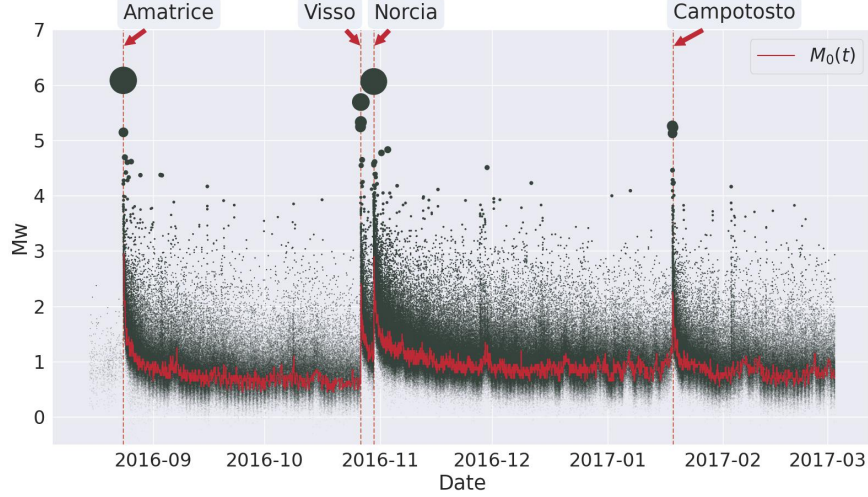


Figure 1: The magnitudes and times of the AVN sequence 2016-2017 (Tan et al., 2021b) used to evaluate the performance of the neural and ETAS model. Marked with a dashed red line are the times of the 4 major events of the sequence. The size of the points are plotted on a log scale corresponding to  $M_w$ . An estimate of the temporal completeness  $M_0(t)$  is plotted using the maximum curvature method (Wiemer & Wyss, 2000).

## 2.1 Amatrice-Visso-Norcia High Resolution Catalog

On the 24th of August 2016 a  $M_w 6.0$  earthquake was recorded near the town of Amatrice in northern Lazio, central Italy. It was followed by a  $M_w 5.9$  near the town of Visso on the 26th of October and a  $M_w 6.5$  near the town of Norcia four days later. Finally, in January 2017, four events between  $M_w 5.0$  and  $M_w 5.5$  struck the Campotosto area. Figure 1 depicts the evolution of this seismic sequence over time.

The INGV produced a routine catalog for the 1 year period containing this sequence (Chiaraluce et al., 2017). Their catalog contains roughly 82,000 events with a completeness of  $Mc = 2.3$  (Mancini et al., 2019). With the use of a neural network based phase picker to determine P and S-wave arrival times, an enhanced catalog has been created for the same earthquake sequence (Tan et al., 2021b). This catalog contains around 900,000 events and has an overall value of completeness of the catalog  $Mc = 0.2$ . We estimated the time varying completeness of this catalog using the maximum curvature method (Wiemer & Wyss, 2000) with samples of 300 events and can see clear variation in completeness particularly following large magnitude earthquakes. This approximate method for the

completeness is only used to show the variability across the catalog and is not used directly in any modelling.

## 2.2 Synthetic Catalog

For the forecasting experiment of synthetic ETAS data, we generate a dataset using the simulator by Mizrahi et al. (2021a), with uniform background intensity  $\mu$  and triggering function,

$$g(t, x, y, m) = \frac{k_0 e^{a(m-M_0)}}{\frac{(t+c)^{1+\omega}}{e^{t/\tau}} ((x^2 + y^2) + d e^{\gamma(m-M_0)})^{1+\rho}}.$$

The ETAS parameters ( $\log_{10} \mu = -6.6$ ,  $\log_{10} k_0 = -3.15$ ,  $a = 2.85$ ,  $\log_{10} c = -2.95$ ,  $\omega = -0.03$ ,  $\log_{10} \tau = 3.99$ ,  $\log_{10} d = -0.35$ ,  $\gamma = 1.22$ ,  $\rho = 0.51$ ,  $M_0 = 1 M_w$ ) are taken close to Mizrahi et al. (2021a) with higher background rate to account for the lower  $M_0$ . The resulting dataset of roughly 250,000 events comes from removing all the spatial co-variates.

We also generate a second synthetic dataset from the first by emulating short-term aftershock incompleteness using the time-dependent formula from Helmstetter et al. (2006),

$$M_0(M, t) = M - 4.5 - 0.75 \log(t),$$

where  $M$  is the mainshock magnitude. Events below the function are removed using the six largest events as mainshocks in this synthetic catalog.

## 3 Theoretical Background

In this section we briefly introduce the theory of neural point processes. We begin with the basic theory of temporal point processes and show how to use this to construct neural point processes for temporal forecasts only. In Section 4, we extend a temporal neural point process to a marked neural process.

### 3.1 Point Processes

A temporal point process (Daley et al., 2003) is a stochastic process that generates a sequence of discrete events at times  $\{t_i\}_{i=1}^n$  in a given observation interval  $[0, T]$ . The process is characterised by its conditional intensity function  $\lambda(t|H_t)$ , which gives the expected number of events in a small interval about  $t$  conditioned on the event history



$$H_t = \{t_i | t_i < t\}:$$

$$\lambda(t|H_t)dt = \mathbb{E}[N([t, t+dt])|H_t].$$

The intensity function (Rasmussen, 2018) completely defines the point process and can take a variety of functional forms. The most basic form is the stationary Poisson process (Daley & Vere-Jones, 2003) which assumes that all events are independent of each other, and the conditional intensity function is constant. Self-exciting point processes assume that events increase the likelihood of subsequent events, with a popular class of these processes being the Hawkes process (Hawkes, 1971). The Hawkes process is defined by its conditional intensity  $\lambda(t|H_t) = \mu + \sum_{t_i < t} g(t - t_i)$ , where  $g(s)$  is a kernel function defining how past events trigger subsequent events.

With the conditional intensity function specified we can write an expression for the log-likelihood of observing a sequence of events  $\{t_i\}_{i=1}^n$ ,

$$\log L(\{t_i\}) = \sum_{i=1}^n \left[ \log \lambda(t_i|H_{t_i}) - \int_{t_{i-1}}^{t_i} \lambda(t|H_t)dt \right]. \quad (1)$$

We assume that the observation interval is  $[t_1, t_n]$  to make the algebra in the remainder more compact. It is trivial to relax this to an arbitrary interval  $[0, T]$ .

Temporal point processes can be extended to incorporate marks. A marked point process is stochastic process that generates events paired with a mark,  $\{t_i, m_i\}_{i=1}^n \in (\mathbb{R}_{>0} \times \mathcal{M})$ . In our setting this represents the occurrence times of earthquakes with their corresponding magnitudes. A marked point process is defined by its conditional intensity function,

$$\lambda(t, m|H_t)dtdm = \mathbb{E}[N(dt \times dm)|H_t],$$

with the log-likelihood of observing a marked sequence of events given by

$$\log L(\{t_i, m_i\}) = \sum_i \left[ \log \lambda(t_i, m_i|H_t) - \int_{t_{i-1}}^{t_i} \int_{\mathcal{M}} \lambda(t, m|H_t)dmdt \right]. \quad (2)$$

### 3.2 Neural Point Processes

The most common form of neural point processes seeks to obtain a compact representation of the event history through the use of an RNN (Du et al., 2016). In this approach, an input representing the inter-event times  $\tau_i = t_i - t_{i-1}$  is first fed into the RNN. A hidden state  $\mathbf{h}_i$  of the RNN is updated

$$\mathbf{h}_i = \sigma(W^h \mathbf{h}_{i-1} + \mathbf{w}^\tau \tau_i + \mathbf{b}^h)$$

where  $\{W^h, \mathbf{w}^\tau, \mathbf{b}^h\}$  are learnable parameters, and  $\sigma$  is an activation function. The conditional intensity function is then formulated as a function of the elapsed time from the

most recent event and is dependent on the hidden state of the RNN,

$$\lambda(t|H_t) = \phi(t - t_i|\mathbf{h}_i),$$

where  $\phi$  is a non-negative function referred to as the hazard function and  $t_i$  is the time of the most recent event. To avoid having to numerically integrate the intensity function directly, Omi et al. (2019) model the integral of the hazard function with a fully connected neural network

$$\Phi(\tau|\mathbf{h}_i) = \int_0^\tau \phi(s|\mathbf{h}_i)ds.$$

With the construction of the model in this way, the log-likelihood of observing a sequence of event times (1) can be written as:

$$\log L(\{t_i\}) = \sum_i \left[ \log \frac{\partial}{\partial \tau} \Phi(\tau_i|\mathbf{h}_i) - \Phi(\tau_i|\mathbf{h}_i) \right].$$

## 4 Methods

Since the neural point process introduced by Omi et al. (2019) is purely temporal, to model seismicity we must extend their model to our requirements. Particularly we require that forecasts be dependent on the history of both times and magnitudes, as magnitudes are an important predictor of seismicity (Utsu, 1970, 1971). We also require a forecast over the next magnitude where, unlike the Gutenberg-Richter law (Gutenberg & Richter, 1936) routinely assumed in the ETAS framework, this is also dependent on the history of events. Finally, we require that we may make forecasts of earthquakes above some target magnitude threshold despite a dependence on earthquakes below that target threshold. In Section 4.1, we first extend the structure of the neural network by Omi et al. (2019) to maximise the likelihood of observing a marked sequence of events, including constructing a time-history dependent magnitude distribution. In Section 4.2, we show how we adjust this new structure to target events above a magnitude threshold.

To aid in the development of more flexible forecasting models, we will make both the dataset and models used in this study available after publication on <https://github.com/ss15859/Neural-Point-Process>.

#### 4.1 Continuously Marked Neural Point Process

We begin with the factorisation of the joint conditional intensity function into its marginal intensity and conditional density function, following Daley et al. (2003),

$$\lambda^*(t, m) = \lambda^*(t)f^*(m|t),$$

where  $\lambda^*(t)$  is the marginal conditional intensity function of  $t$ , and  $f^*(m|t)$  is the conditional density function of the mark at time  $t$ . Both of these functions are dependent on the history  $H_t$ , here denoted by the asterisk \*. To construct the likelihood for the marked sequence we model these two functions separately.

With the factorisation, the expression for the log-likelihood of observing the marked sequence of events (2) becomes,

$$\log L(\{t_i, m_i\}) = \sum_i \left[ \log \lambda^*(t_i) + \log f^*(m_i|t_i) - \int_{t_{i-1}}^{t_i} \lambda^*(t) dt \right]. \quad (3)$$

Now with a two dimensional input, the hidden state of the RNN is updated as a linear combination of the inter-event times and magnitudes. This is the continuous mark extension to the RNN update from Du et al. (2016),

$$\mathbf{h}_i = \sigma(W^h \mathbf{h}_{i-1} + \mathbf{w}^\tau \tau_{i-1} + \mathbf{w}^m m_{i-1} + \mathbf{b}^h),$$

where  $\mathbf{w}^m$  is an additional learnable parameter.

The marginal intensity function is formulated as a function of the elapsed time from the most recent event and is dependent on the hidden state of the RNN (Du et al., 2016),

$$\lambda(t|H_t) = \phi(\tau = t - t_i|\mathbf{h}_i), \quad (4)$$

where  $\phi$  is a non-negative function referred to as the hazard function and  $t_i$  is the time of the most recent event. Following Omi et al. (2019), we model the cumulative hazard function using a fully connected neural network,

$$\Phi(\tau|\mathbf{h}_i) = \int_0^\tau \phi(s|\mathbf{h}_i) ds.$$

which allows us to differentiate this with respect to  $\tau$  to extract the hazard function. The derivative is easily obtained through automatic differentiation (Van Merriënboer et al., 2018), which is available in all neural network packages.

We now also formulate the conditional density function of the mark at time  $t$  as a function of the current mark. This is dependent on the time since the most recent event and the hidden state of the RNN,

$$f(m|t, H_{t_i}) = \psi(m|\tau, \mathbf{h}_i).$$

We again model its cumulative distribution with a fully connected neural network,

$$\Psi(m|\tau, \mathbf{h}_i) = \int_0^m \psi(\mu|\tau, \mathbf{h}_i) d\mu.$$

Although this integral does not feature in the expression for the log-likelihood, we still opt for this approach over directly modelling the density function with a neural network. This follows from work on neural density estimation where positive weights can be enforced in the network to capture the positivity and monotonicity of cumulative distribution functions (Chilinski & Silva, 2020). We can then obtain the density function again through automatic differentiation.

We can now write the log-likelihood as:

$$\log L(\{t_i, m_i\}) = \sum_i \left[ \log \lambda^*(t_i) + \log f^*(m_i|t_i) - \int_{t_{i-1}}^{t_i} \lambda^*(t) dt \right] \quad (5)$$

$$= \sum_i \left[ \log \phi(\tau_i|\mathbf{h}_i) + \log \psi(m_i|\tau_i, \mathbf{h}_i) - \int_0^{t_i - t_{i-1}} \phi(t|\mathbf{h}_i) dt \right] \quad (6)$$

$$= \sum_i \left[ \log \frac{\partial}{\partial \tau} \Phi(\tau_i|\mathbf{h}_i) + \log \frac{\partial}{\partial m} \Psi(m_i|\tau_i, \mathbf{h}_i) - \Phi(\tau_i|\mathbf{h}_i) \right]. \quad (7)$$

We model both the cumulative hazard function  $\Phi$ , and the conditional distribution function of the marks  $\Psi$ , using a feed-forward neural network. The network depicted in Figure 2 consists of four component parts. The first part is the recurrent section, which finds an encoding of the history of the point process. The output of the recurrent section  $\mathbf{h}_i$  passes into two fully connected components. One models the integral of the intensity function, this is a function of the time of the next event  $\tau$ . The other component models the integral of the conditional density function of the next mark, this is dependent on the next time  $\tau$ , but is a function of the next mark  $m$ .

Since passing long sequences into recurrent neural networks can often lead to exploding or vanishing gradient problems (Hochreiter, 1998), we do not pass the whole history of the point process into the recurrent section, but the past  $d$  events. This implies that we use only the past  $d$  events to forecast the next event. Thus, this model is learning to estimate the intensity given a recent history of  $d$  events. This hyperparameter is kept the same as Omi et al. (2019) at  $d = 20$ . A naive tuning search found no significant improvement at larger values of  $\{50, 100, 200, 500\}$ . This difference to the full-history ETAS model is discussed in Section 6.1.

We enforce positive weights in both fully connected sections of the network to capture the positivity and monotonicity that is required by both cumulative functions. In the final component of the network we formulate the output as the log-likelihood of ob-

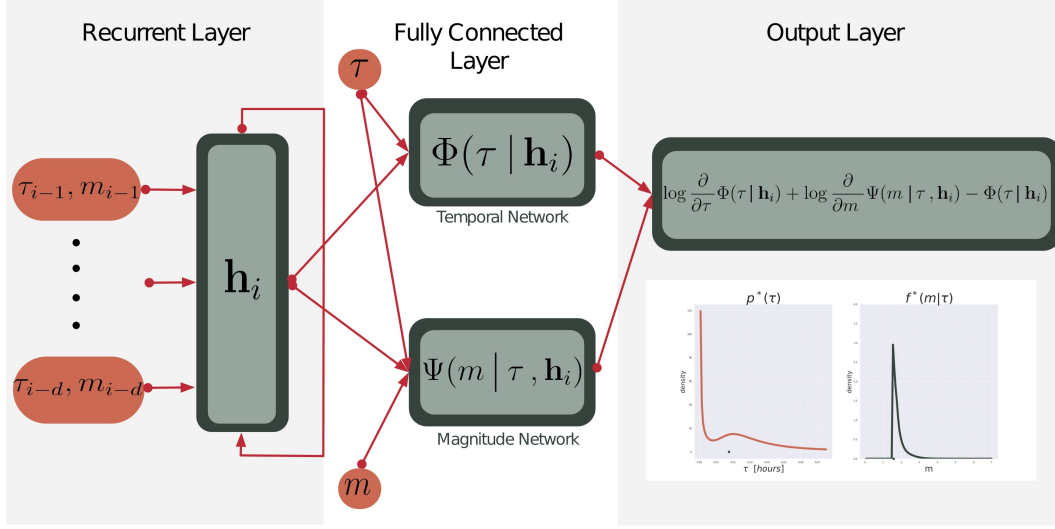


Figure 2: The proposed network comprises four sections. First, the inter-event times and magnitudes of the last  $d$  events are fed into a recurrent section consisting of 64 recurrent units. The output of this section is fed into two fully connected sections where it is combined with the next inter-event time  $\tau$  for the temporal network and additionally with the next magnitude  $m$  for the magnitude network. The outputs of both these sections are combined to formulate the log-likelihood of the next inter-event time and magnitude  $\{\tau, m\}$ . We can separate the temporal and magnitude terms in this likelihood to give point evaluations of the density of the next inter-event time and conditional density of the next magnitude. The dependence structure of the point process is expressed by the connections between sections in the network.

serving the pair  $\{\tau, m\}$ . To construct this output, one backward pass is performed to calculate the derivatives with respect to the next time and the next magnitude found in eq (7). This output is exactly what is maximised to learn the weight parameters, during a second backwards pass.

## 4.2 Target Events

The growth in machine learning generated catalogs from their predecessors is found through detecting events of magnitude much lower than previously possible. However, operationally, we may not care about the forecasting of these smaller events if it is the larger ones that are the most hazardous. The aim is therefore to use the smaller events to forecast earthquakes above some target threshold. In this section we outline how this

is done for both ETAS and the neural model. We hereafter call events above the target magnitude threshold  $M_d$  target events.

Let  $\{(t_i, m_i)\}_{i=1}^n \in [0, T] \times \mathcal{M}$  be the entire sequence of events, complete down to  $M_0$ . We seek to make forecasts of events above magnitude  $M_d$ . This corresponds to the sequence:

$$\{(t_j, m_j) : m_j \geq M_d\}_{j=1}^k.$$

To ease in the distinction between ‘all events’ and the target events, we subscript the former with  $i$  and the latter with  $j$ .

Let  $\lambda_0(t, m|H_t)$  denote the joint intensity function of all events above  $M_0$ . We seek to learn the intensity of events above  $M_d$ , denoted  $\lambda_d(t, m|H_t)$ , where the history  $H_t$  contains all events  $\{(t_i, m_i)\}_{i:t_i < t}$ . The ground intensity above the target threshold is found by marginalising the joint intensity over the target magnitude region,

$$\lambda_d(t|H_t) = \int_{M_d}^{\infty} \lambda_0(t, m|H_t) dm.$$

The log likelihood of target events is then given by:

$$\log L(\{t_j, m_j\}) = \sum_{j:m_j \geq M_d} \left( \log \lambda_d(t_j, m_j|H_{t_j}) - \int_{t_{j-1}}^{t_j} \lambda_d(t|H_t) dt \right).$$

For ETAS, the rate above magnitude  $M_d$  is a fraction of the rate above  $M_0$ , due to the independent distribution for magnitudes,

$$\lambda_d(t|H_t) = \int_{M_d}^{\infty} \lambda_0(t, m) dm = \left( \int_{M_d}^{\infty} f_{GR}(m) dm \right) \lambda_0(t|H_t) = p_d \cdot \lambda_0(t|H_t),$$

where  $f_{GR}(m)$  is the Gutenberg-Richter law and  $p_d$  is simply the probability that  $m \geq M_d$ . Therefore the expression for the likelihood is relatively simple,

$$\log L(\{t_j, m_j\}) = \sum_{j:m_j \geq M_d} \left[ \log (p_d \cdot \lambda_0(t_j, m_j|H_{t_j})) - \int_{t_{j-1}}^{t_j} p_d \cdot \lambda_0(t|H_t) dt \right].$$

For the neural model, we make use of the fact that the integral of the intensity function between target events,  $\{(t_j, m_j) : m_j \geq M_d\}_{j=1}^k$ , can be expressed as a sum of dis-

joint integrals between all events  $\{(t_i, m_i)\}_{i=1}^n$ ,

$$\log L(\{t_j, m_j\}) = \sum_{j:m_j \geq M_d} \left[ \log \lambda_d(t_j, m_j | H_{t_j}) - \int_{t_{j-1}}^{t_j} \lambda_d(t | H_t) dt \right] \quad (8)$$

$$= \sum_{j:m_j \geq M_d} [\log \lambda_d(t_j | H_{t_j}) + \log f_d(m_j | t_j, H_{t_j})] - \int_{t_0}^{t_k} \lambda_d(t | H_t) dt \quad (9)$$

$$= \sum_{\substack{i:m_i \geq M_0 \\ t_i \leq t_k}} \left[ (\log \lambda_d(t_i | H_{t_i}) + \log f_d(m_i | t_i, H_{t_i})) \mathbf{I}\{m_i \geq M_d\} - \int_{t_{i-1}}^{t_i} \lambda_d(t | H_t) dt \right] \quad (10)$$

$$= \sum_{\substack{i:m_i \geq M_0 \\ t_i \leq t_k}} \left[ \left( \log \frac{\partial}{\partial \tau} \Phi(\tau_i | \mathbf{h}_i) + \log \frac{\partial}{\partial m} \Psi(m_i | \tau_i, \mathbf{h}_i) \right) \mathbf{I}\{m_i \geq M_d\} - \Phi(\tau_i | \mathbf{h}_i) \right], \quad (11)$$

where between (8) and (9) we have factorised the joint intensity of events above  $M_d$ ,  $\lambda_d(t, m | H_t)$ , into the ground intensity above the target threshold,  $\lambda_d(t | H_t)$ , and the distribution of the next magnitude above the target threshold given the time and the history,  $f_d(m | t, H_t)$ . Between (9) and (10) we changed the summation from being over target events to being over all events by adding the indicator function  $\mathbf{I}\{m_i \geq M_d\}$ . The integral in (9) becomes the sum of integrals in (10) through a decomposition into disjoint integrals between all events  $\{(t_i, m_i)\}_{i=1}^n$ . Thus the neural model may target events by adding the indicator function  $\mathbf{I}\{m_i \geq M_d\}$  to the expression for the log-likelihood. Now the hazard function models the rate above  $M_d$  as a function of the time from the last event (of any magnitude),

$$\lambda_d(t | H_t) = \phi(\tau = t - t_i | \mathbf{h}_i).$$

### 4.3 Experimental Design

For both the synthetic data and real data we apply the same training and testing procedure illustrated in Figure 3. At a fixed point in time along the sequence we set a marker and train on data up to that point. Following that, the remainder of the sequence will be used as the test set. For the synthetic catalogs this is done at one single point in time, whereas for the Central Apennines sequence we make three partitions - each just before the Visso, Norcia and Campotosto earthquakes. By making these partitions we can see how the performance of each model is affected by the number of training datapoints as well as the number of major earthquakes.

We seek to understand how different magnitude thresholding affects the performance of each of the models. For each of the partitions, we look at the performance of both mod-

els as the magnitude threshold of the input catalog is lowered, a parameter we refer to as  $M_{\text{cut}}$ . For the AVN catalog and incomplete synthetic catalog, this includes lowering  $M_{\text{cut}}$  into periods where  $M_{\text{cut}} < M_0(t)$ . We keep fixed the magnitude of events we wish to target at  $M_d = 3$  Mw.

Both models are trained by maximum likelihood estimation (MLE) on the training dataset. For ETAS we use the intensity function defined by Ogata (1988) and maximise the likelihood through Nelder-Mead optimisation, chosen for its robustness. For the neural model, we maximise the likelihood defined in equation (11) through ADAM optimisation (Kingma & Ba, 2014) written in Tensorflow (Abadi et al., 2015).

We compare the two models' performance through the log-likelihood of the events in the testing set. We separate the temporal and magnitude terms in the likelihood equation (3), to analyse their predictive skills on the two target variables separately. To compare the performance across different magnitude thresholds, we also fit a homogeneous Poisson model. For the temporal log-likelihood we can therefore present the log-likelihood gain from a benchmark Poisson model, whereas for the magnitude forecasts, simply the log-likelihood is reported. We shall refer to both performance metrics as log-likelihood scores and to make general comparisons across the models, we compare the mean log-likelihood score per earthquake as well as construct a 95% confidence interval to assess the variability. The confidence interval is constructed with 1000 bootstrap samples of the log-likelihood scores.

## 5 Results

### 5.1 Synthetic Data

Despite the synthetic catalog being complete down to  $M_0 = 1$  Mw, we only lower  $M_{\text{cut}}$  down to 1.7 due to the computational time it takes to find the MLE parameters of ETAS for such a large dataset. Figure 4 shows the computation time (CPU hours) to train each of the models as a function of the size of the training set using an 2.4 GHz Intel E5-2680 v4 (Broadwell) CPU. The neural model is significantly faster to train than ETAS due to the likelihood function not being dependent on the full history of the sequence, giving it complexity  $O(N)$  (Shchur et al., 2021). ETAS in contrast has complexity  $O(N^2)$  due to the double sum in the likelihood.

Figure 5 shows 95% confidence intervals for the log-likelihood scores on the synthetic catalogs for the varying magnitude thresholds. By varying the value of  $M_{\text{cut}}$  the



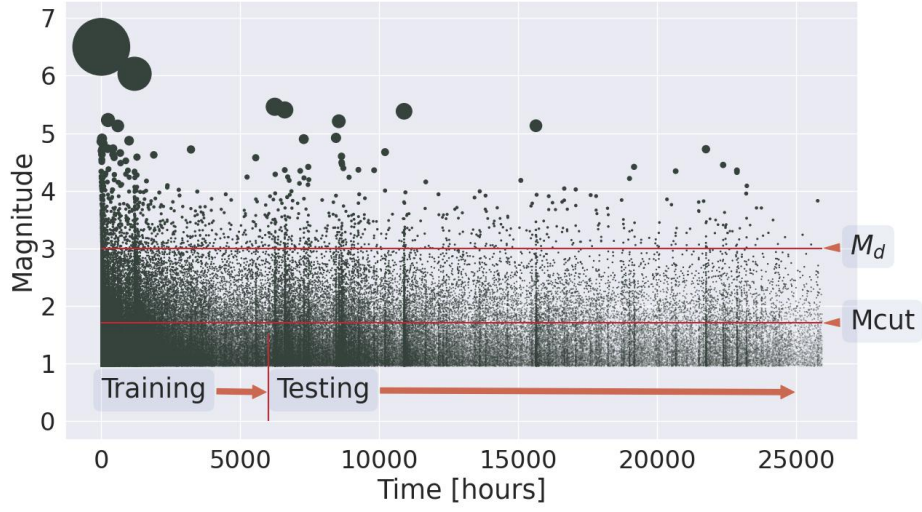


Figure 3: The synthetic catalog with an outline of the training and testing procedure. We train up to a fixed point in time in the catalog, following which the remainder of the catalog is used for testing. We vary the value of the threshold for the input catalog ( $M_{\text{cut}}$ ) and keep fixed the value of the target threshold ( $M_d$ ).

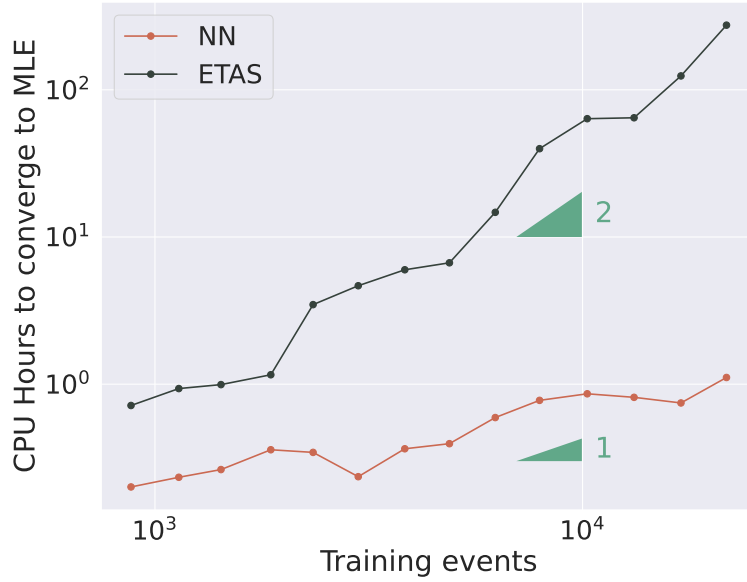


Figure 4: Time on a single CPU required to train each of the models as the training size increases. Each model is trained by maximising the likelihood of the training data.

size of the training dataset changes, depicted by the green barplots. In Figure 5a) the temporal log-likelihood scores for both models on the complete synthetic catalog are displayed. Although there are fluctuations, there is no significant difference between the two models' log-likelihood gain from the same Poisson baseline for all values of  $M_{\text{cut}}$ , suggesting the neural model has learnt to capture the ETAS data sufficiently well. Although the mean of ETAS increases as we lower  $M_{\text{cut}}$ , whereas the mean of the neural model fluctuates, these changes are non-significant. We speculate that we do not see significant improvement as we lower  $M_{\text{cut}}$  due to the fact that we are fitting temporal models to spatio-temporal data.

Figure 5b) shows the temporal log-likelihood scores for the incomplete synthetic catalog. Just as in Figure 5a), ETAS remains constant down to  $M_{\text{cut}}$  2.0. But now, on this incomplete dataset, the performance of ETAS significantly reduces as  $M_{\text{cut}}$  is lowered below this threshold. In contrast, the neural model remains constant in performance and significantly outperforms ETAS for all but the highest two thresholds, demonstrating a robustness to the missing data in this catalog. By synthetically recreating incompleteness we remove many datapoints, therefore we can fit ETAS to a lower  $M_{\text{cut}}$  as we do not experience the longer training times of the complete catalog.

Figure 5c) shows the magnitude log-likelihood scores for the complete synthetic catalog. For the highest two thresholds, the neural model performs significantly worse than ETAS but then remains marginally worse for all other values of  $M_{\text{cut}}$ . For the magnitude scores for the incomplete data in Figure 5d), ETAS significantly outperforms the neural model at the higher thresholds. As the threshold is lowered the two perform more similarly, owing to an increase in performance from the neural model and a slight decrease from ETAS.

We can understand the marginal difference in performance between the two models at lower thresholds by looking at their respective distributions. Figure 6 shows five instances of the magnitude distribution learnt by both models at  $M_{\text{cut}} = 1.7$  for the complete catalog. For ETAS we simply learn the  $b$  value of the Gutenberg-Richter (GR) law whereas for the neural model, a history and time-dependent distribution for the next magnitude is learnt. In these five instances, although the neural model can closely approximate the GR law, allowing it to be time-history dependent means that its predictions vary across different occurrences in the sequence and therefore in this synthetic example performs worse than the stationary data-generating distribution. Since the neural

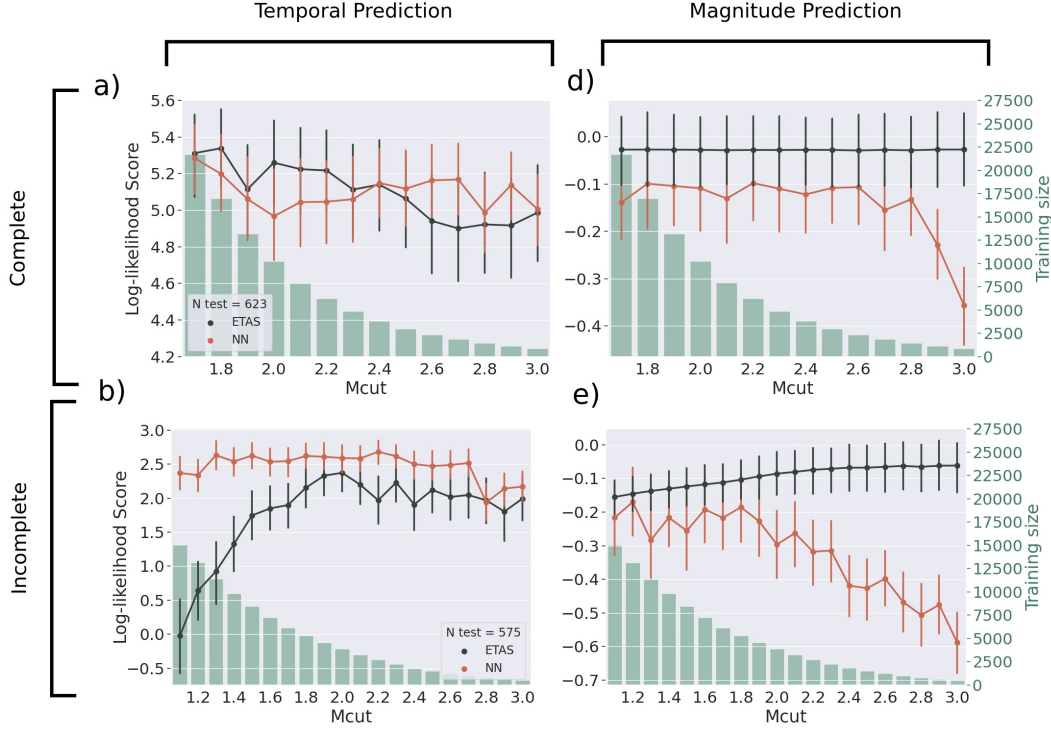


Figure 5: Results from the synthetic tests. 95 % confidence intervals for the log-likelihood scores for each model as a function of  $M_{\text{cut}}$  (the magnitude threshold of the input catalog). The size of the training set is displayed in the green barplot; the size of the testing set in the legend. a) temporal log-likelihood gain from Poisson for the complete synthetic catalog. b) temporal log-likelihood gain for the incomplete catalog. c) magnitude log-likelihood for the complete catalog. d) magnitude log-likelihood for the incomplete catalog.

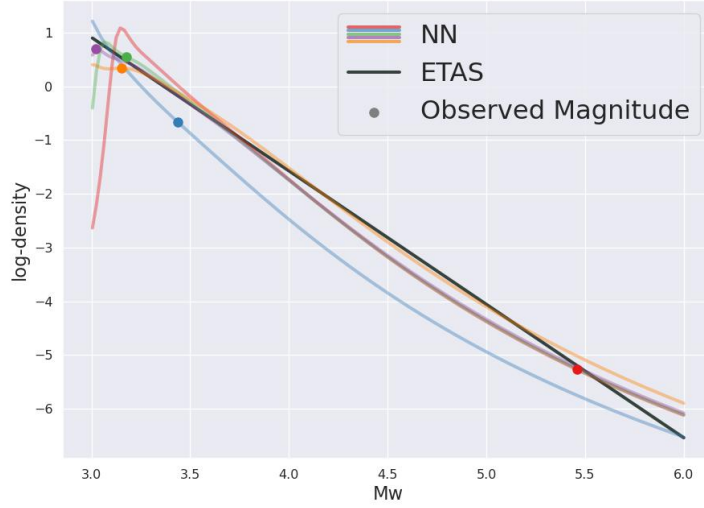


Figure 6: Five examples of the forecasted magnitude distributions from the complete synthetic catalog tests at  $M_{\text{cut}} = 1.7$  compared with the ETAS Gutenberg-Richter law. The magnitudes of the observed events are plotted as points along the log-density for the neural model.

model contains orders of magnitude more parameters, this result indicates overfitting (Lawrence et al., 1997).

## 5.2 AVN Catalog

Figure 7 shows the log-likelihood scores for both models on each of the testing-training partitions on the AVN catalog, where Figure 7a) is for training both models up to the Visso 5.9 Mw event, Figure 7b) up to the Norcia 6.5 Mw event and Figure 7c) up to the first of the major Campotosto earthquakes. Across all training-testing partitions Figure 7a)-7c), as  $M_{\text{cut}}$  is lowered below 3 Mw, the performance of ETAS decreases consistently. The neural model, however, either remains constant in performance or improves as  $M_{\text{cut}}$  is lowered. In addition, as the neural model is trained on a longer period of time, its performance improves. For higher values of  $M_{\text{cut}}$  the neural model performs significantly worse than ETAS when trained up to Visso, but with the additional training data leading up to Norcia and Campotosto, it is similar to ETAS. For low values of  $M_{\text{cut}}$ , the performance of the neural model is significantly better than ETAS. When comparing across

all values of  $M_{\text{cut}}$  neither model is significantly better than the other. Generally, the neural model is more robust to different values of  $M_{\text{cut}}$  than ETAS. Figure S1 of the Supporting Information shows how the fitted ETAS parameters change with  $M_{\text{cut}}$ .

The magnitude log-likelihood scores in Figure 7d)-7f) show that the time-history dependent magnitude distribution generally cannot match the predictive power of the stationary GR law. The performance of ETAS remains constant for all values of  $M_{\text{cut}}$  and testing-training partitions. In Figure 7d) and Figure 7e) the neural model improves in performance as  $M_{\text{cut}}$  is lowered, where it only performs closely to ETAS at the very lowest threshold. This and the fact that it performs much closer to ETAS when training up to Campotosto (Figure 7f)) suggests that it is learning and improving when shown more data.

To compare the models' performance as a function of time, Figure 8 displays the cumulative information gain (CIG) of the neural model over ETAS, for both models trained up to the Norcia earthquake. This information gain is simply the difference in the log-likelihood scores, where we subtract the score of ETAS from the neural model for both the magnitude and event-time term of the likelihood. The CIG is plotted per earthquake, but the evolution with time since the Norcia earthquake is also displayed. Figure 8a) shows the CIG for event time forecasts. Beyond the trend that the neural model improves over ETAS as we lower  $M_{\text{cut}}$ , the improvement varies over the testing catalog. For the thresholds that give the largest gain, ( $M_{\text{cut}} = 1.2, 1.4, 1.6, 1.8$ ), the period of time with the greatest amount of gain, indicated by the steepest gradient of the curve, is found within the first 2 hours of the Norcia earthquake. This is followed by a reduced improvement up to 24 hours, beyond which it levels out and remains relatively linear.

Figure 8b) shows the CIG for the magnitude forecasts, confirming the loss in average performance of the neural model over ETAS. All thresholds decrease fairly steadily for nearly all of the testing period, apart from immediately following the Norcia earthquake. For the lower thresholds the period of time following Norcia sees an improvement over ETAS very briefly before declining.

Figure 8c) shows the IG of the neural model over ETAS but now as a function of the estimate of the completeness of the testing period. Both models are trained up to Norcia for  $M_{\text{cut}} = 1.2, 2.0, 2.8$ . A locally weighted scatterplot smoothing (lowess) regression (Cleveland, 1979) with 95 % confidence intervals estimates this relationship. For  $M_{\text{cut}} = 1.2$ , the difference between the two models is smallest for intermediate values

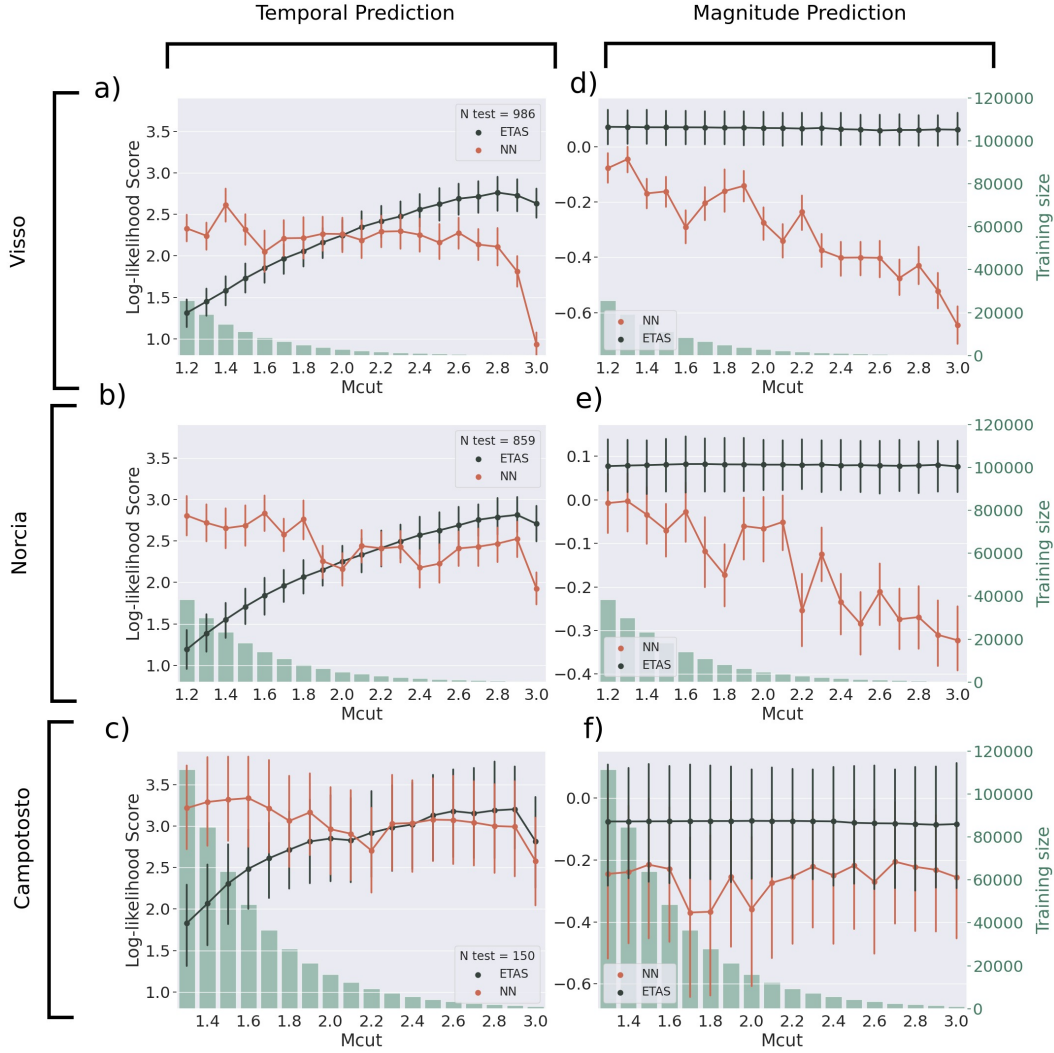


Figure 7: Results from tests on AVN catalog. 95 % confidence intervals for the log-likelihood score of each model for varying values of Mcut. The size of the training set is displayed in the green barplot as well as the size of the testing set in the legend. a)-c) depicts the temporal log-likelihood gain from Poisson. In a), both models are trained up to the Visso earthquake, in b) both models are trained up to the Norcia earthquake and in c) both are trained up to the Campotosto earthquakes. d) - f) depict the magnitude term of the log-likelihood for the same training-testing partitions.

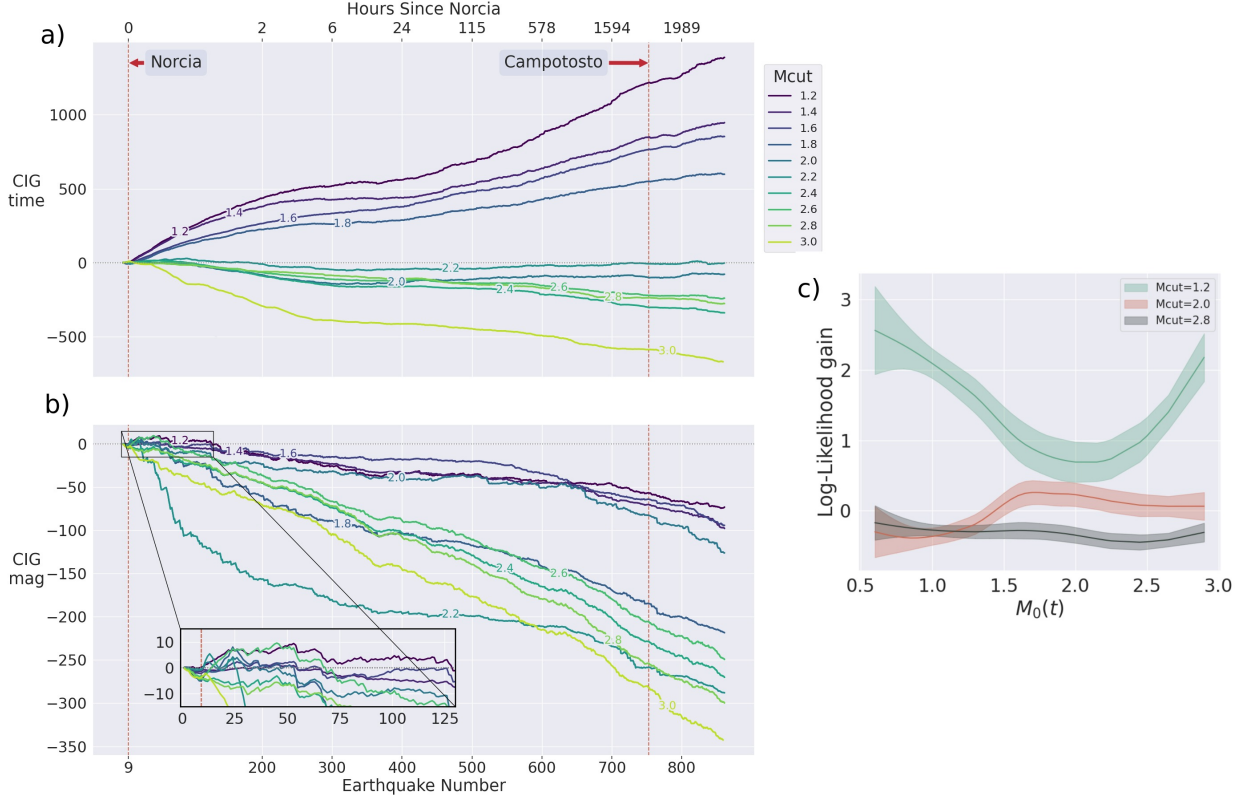


Figure 8: a) - b) The Cumulative Information Gain (CIG) of the neural model over ETAS for a range of values of  $M_{cut}$ . The models are trained up to the Norcia earthquake and the plot depicts the evolution of the CIG from the Norcia earthquake to the end of the catalog. The curve is plotted per event, however, the actual time since the Norcia earthquake is displayed on the top axis. a) displays the CIG for event-time forecasts, b) displays the CIG for magnitude forecasts. c) displays the information gain of the neural model over ETAS as a function of the completeness of the testing catalog - both models are trained up to Norcia for  $M_{cut} = 1.2, 2.0, 2.8$ .

of the incompleteness (around  $M_0(t) = 2.0$ ), but for the most complete ( $M_0(t) = 1.0$ ) and most incomplete ( $M_0(t) = 3.0$ ) parts of the testing catalog, the neural model performs greatest compared to ETAS. At  $M_0(t) = 2.0$  where the relative performance of ETAS is best, the confidence interval for the log-likelihood difference between the two models lies above zero, centered around 0.75. So, there is no value of completeness in the testing catalog where ETAS performs as well as the neural model. For  $M_{cut} = 2.0$ ,

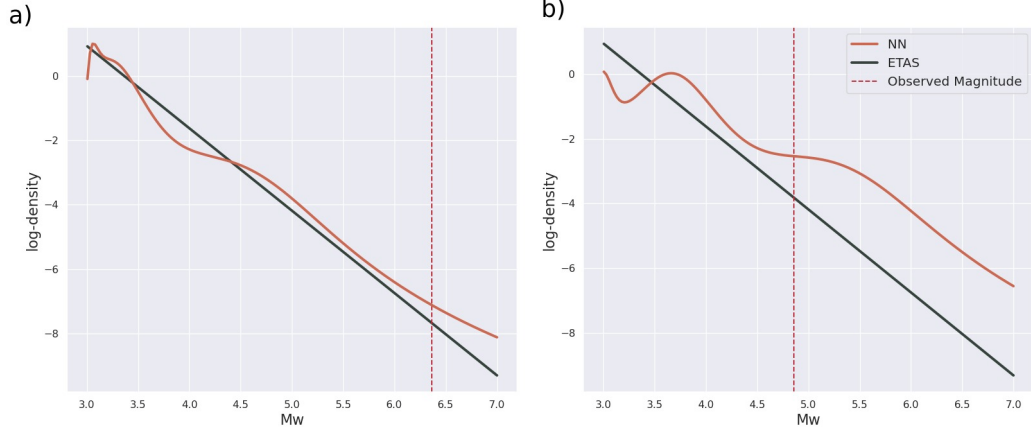


Figure 9: The forecasted magnitude distribution of each model, a) at the occurrence of the Norcia earthquake, and b) at the occurrence of the next Mw3+ earthquake following Norcia.

the neural model outperforms ETAS during more incomplete periods of the testing period and for  $M_{\text{cut}} = 2.8$ , ETAS is consistently better across all values of completeness.

The forecasted magnitude distribution of each model shown in Figure 9 provides some insight into the cause of this immediate improvement over ETAS’s GR forecast right after Norcia. Figure 9a) shows the magnitude distribution of each model at the time of the Norcia earthquake. The two distributions resemble each other relatively closely. At the time of the next Mw3+ event, Figure 9b), the neural model has shifted its density towards higher magnitude values anticipating a lack of observed smaller magnitude earthquakes due to catalog incompleteness.

## 6 Discussion

### 6.1 Approximating ETAS

The ability to approximate ETAS data using a neural point process is a benchmark goal. Demonstrating a baseline level of expressiveness is essential before any work on real data is done. Given other neural point process models regularly use univariate or discretely marked Hawkes data as a baseline (Du et al., 2016; Omi et al., 2019), this result provides an example of fit to continuous bivariate Hawkes data.



Specifically, the merit of this fit to ETAS data is that it uses a truncated history of events. Truncating the ETAS model to sum over only the last  $d$  events would dramatically change the evaluation of the intensity function between a significantly large earthquake  $d$  events ago or  $d+1$  events ago, thus instead the way in which to formulate the relationship between the intensity and these truncated events is learnt. The past  $d$  events are exhibiting behaviour based on the events prior and thus we do not directly specify the contribution from events further back in time, instead dependence on such events is learnt indirectly.

The reduction in the amount of history each forecast is dependent upon drastically improves the computational requirements for both learning and evaluation of the likelihood. To create forecasting models alongside the growing size of earthquake catalogs generated through machine learning based phase picking, we require models that scale well with the data. Shown here is that we can achieve the same predictive accuracy as ETAS (in a synthetic catalog) but with a smaller computational budget.

## 6.2 Embracing and Ignoring Data Incompleteness

At the previously explored thresholds of this sequence, (Ebrahimian & Jalayer, 2017; Marzocchi et al., 2017; Mancini et al., 2019, 2022), the neural point process performs similarly to ETAS. The biggest deviations between the two models are found as the magnitude threshold is lowered into new unexplored regions revealed by this enhanced catalog. The deviations come from the fact that the neural model increases gradually in performance as the threshold is lowered whereas ETAS drastically decreases in performance. Below, we offer an interpretation for these results.

We argue that the largest gains made by the neural model are due to its ability to handle the incomplete data immediately following large earthquakes. There are two justifications for this: Given that the magnitude threshold of the input catalog is lowered into regions where there are periods when  $M_0(t) > M_{\text{cut}}$ , we expect there to be biases in fitting ETAS (Zhuang et al., 2017; Seif et al., 2017; Hainzl, 2016b). The consequences of these biases on ETAS are reflected in the log-likelihood scores on the synthetic incomplete data figure 5b). In this synthetic catalog with short-term aftershock incompleteness, ETAS drastically reduces in performance in contrast to the neural model. Since a basic ETAS is formulated as completely observing all potentially triggering earthquakes it poorly captures incomplete sequences. The same shape of plot is found in the real data

in Figure 7a)-c), where the performance of ETAS decreases as  $M_{\text{cut}}$  is lowered. In contrast, other studies have found decreasing the minimum triggering magnitude improves the performance of ETAS (e.g., Werner et al. (2011); Helmstetter and Werner (2014)). A consequence of the bias in the ETAS parameters is that the forecasting performance is only competitive with the neural model during intermediate values of incompleteness, Figure 8c). Even during complete periods in the testing catalog, since ETAS has been fit on incomplete data, it fails to forecast well.

The second justification comes through considering the process by which the observed data are generated, e.g. as described by Omi et al. (2014). The relationship between the underlying process  $\lambda(t|H_t)$  and the observed process  $v(t|H_t)$  that forms the catalog itself can be written as

$$v(t|H_t) = \lambda(t|H_t) r(t|M_0),$$

where  $r(t|M_0)$  is the probability of detection at time  $t$ . For ETAS variants that deal with time-varying completeness (e.g. Omi et al., 2014; Hainzl, 2016a, 2016b; Mizrahi et al., 2021b), this function has to be estimated alongside the parameters of ETAS. When fitting a temporal ETAS model to data with temporal incompleteness, bias in the fitted parameters comes from the modeling assumption that  $v^*(t) = \lambda^*(t)$ . Through the use of a flexible model such as this neural point process, rather than trying to learn both the underlying process and the detection rate, we instead directly learn the observed process. So in fact, in the construction of the model rather than equation (4), the observation process is approximated,

$$\phi(t - t_i | \mathbf{h}_i) = v^*(t),$$

where  $t_i$  is the time of the last event. As we lower  $M_{\text{cut}}$  into regions of temporal incompleteness, unlike ETAS, the performance of the neural model does not decrease, Figure 5b). This demonstrates the neural models' ability to fit to observed data as it is not biased by an increasing amount of missingness.

Learning to model the observation process requires the assumption that the process of the incompleteness is stationary for future forecasts, thus if there is new methodology in data collection, the model would have to be re-trained on this new data. This is similar to detection rate based methods, (Omi et al., 2014; Hainzl, 2016a, 2016b; Mizrahi et al., 2021b), that would also need to update their detection function with new observational methodology.

To further test the effectiveness of the neural model, a comparison with other ETAS models that specifically deal with incompleteness is needed. But given that methods that deal with incompleteness only increase the computational requirements upon fitting a basic ETAS model, neural point processes could offer a more efficient way to deal with missing data. This is especially important when moving towards using enhanced catalogs such as the one used here. Temporal variations in completeness must be considered when using these catalogs and to be able to use these catalogs we must take more care with the computational efficiency of models.

Although data incompleteness is the most reasonable argument for the significant gains of the neural point process at low magnitude thresholds, we shouldn't limit ourselves to this description. We have reached this conclusion by extrapolating the forecasting results on synthetic data and through arguments about which statistical process is being approximated by the neural model. However, we shouldn't rule out the possibility that the new low magnitude data in this enhanced catalog has contributed additional signal that is not explainable by ETAS. Even for the intermediate value of completeness that results in the best performance of ETAS compared with the neural model, there is still an average log-likelihood gain of around 0.75, Figure 8c). We can compare this with Figure S3c) in the supporting information, where on the incomplete synthetic data we see a truncated curve of the same shape. In this synthetic experiment, there are periods of completeness where the performance between the two models is comparable, however on the real data, since there is no value of completeness that results in comparable performance, this would suggest that something additional is contributing to the gains. The question of whether there is additional signal in low magnitude events found in enhanced catalogs such as this one needs further attention beyond this study. We believe that further development in neural point processes will aid modellers in analysing this wealth of new data as neural models provide more flexible modelling alternatives and can cope with the scale of new enhanced catalogs.

### 6.3 Limitations

This study presents a flexible model that does not suffer from the same misspecification as ETAS due to short-term aftershock incompleteness. Although the size of the gains for these magnitude thresholds is large, we found, however, no significant overall improvement in forecasting ability over ETAS across magnitude thresholds. Comparing

the value of  $M_{\text{cut}}$  that gives the greatest performance for each model finds that although the mean of the neural model is highest, the gain over ETAS is not significant. In this two dimensional time-magnitude domain, given the flexibility of this neural network and the data volume provided, there is insufficient signal in the data to learn anything significantly better than ETAS. This would suggest that the time and magnitude data from low magnitude events alone does not give us additional information in forecasting M3+ events. This motivates considering whether additional features can aid in the forecasting ability of neural point processes. Given that operationally we also require spatial forecasts, this is an obvious future extension to the model. It is natural to expect that including spatial covariates would improve forecasting performance (Utsu, 1955; Ogata, 1998), however, it is not clear that considering them as an additional dimension to the input of an RNN would learn any spatial structure from the data. Neural point processes for spatio-temporal data do not utilise RNNs which are primarily sequence encoders and instead consider models based on Ordinary Differential Equations (ODEs) (Chen et al., 2020; Biloš et al., 2021). We believe such models should outperform RNN-based models on spatio-temporal data.

By modelling the magnitudes by a completely unconstrained density that is also time-history dependent, we create lots of potential for over-fitting to the data (Ying, 2019). This is exactly what is observed during the tests on synthetic data where we learn a ‘noisy’ Gutenberg-Richter law. It is the likely source for the performance which is on average worse than ETAS on the AVN catalog. However, by letting this function be unconstrained, the model was able to make improvements over ETAS immediately following Norcia. This isn’t too surprising since deviations from a stationary GR law have been observed, either through fluctuations of the  $b$ -value in space or time (Schorlemmer et al., 2005; Gulia et al., 2018) or short-term aftershock incompleteness (Kagan, 2004; Woessner & Wiemer, 2005; Helmstetter et al., 2006).

A final limitation of the model presented here is its (in)ability to simulate events into the future. Where simulation of ETAS can be done due its equivalent branching process formulation, simulation from the neural model can only be leveraged through the intensity function at a given point in time with a thinning algorithm (Ogata, 1981). But given that calculation of the intensity with a neural point process is much more efficient than ETAS, it would be quicker to do thinning simulation with a the neural model. Slower simulation can be overcome by using an alternative flexible formulation of the likelihood

such as the one used in Shchur et al. (2019), however, we would lack the ability to target earthquakes above a magnitude threshold.

## 7 Conclusion

We present an initial investigation into the viability of neural point processes for the forecasting of short term seismicity. The neural point process is formulated in a similar way to the ETAS model, only with a much more flexible way of representing the intensity function. Now with much larger earthquake catalogs, data-driven point process models present us with an opportunity to investigate whether these new data may offer some deviation to the parameterization of ETAS as well as providing more computationally efficient models that are robust to missing data.

We extend the existing point process model of Omi et al. (2019) so that it also models the magnitudes associated with the events contained in earthquake sequences. We also show how this model can be used to forecast earthquakes above some target threshold magnitude through decomposing the cumulative hazard function between target events. A notable feature of the presented model is that a forecast is only dependent on a fixed length vector representing the history of events, making the evaluation of the likelihood scale linear with the sample size.

With an experiment on data simulated from the ETAS model we demonstrate this computational advantage by showing a stark improvement on the time to train the neural point process against the ETAS model, whilst still obtaining a similar likelihood score on test data. We find that defining a more flexible time-history dependent magnitude distribution leads to overfitting and consequentially the magnitude likelihood scores are worse than when using a stationary Gutenberg-Richter law.

Through artificially removing events from the synthetic catalog we create a dataset that mimics short-term aftershock incompleteness. We find that on this altered catalog, the performance of ETAS now decreases as the magnitude threshold is lowered. In contrast, the neural model remains constant in performance, suggesting that it is more robust to the missing data found in typical earthquake catalogs.

On real data from the Amatrice-Visso-Norcia sequence the performance of both models vary with respect to the magnitude threshold of the input catalog. Both models perform similarly at previously explored thresholds ( $M_w3+$ ), but when lowered into magnitude regions revealed by the new catalog, ETAS decreases in performance unlike

the neural point process. We argue that these gains are due to the neural model’s ability to handle the incomplete data found in this enhanced catalog. This experiment both motivates the need for considering temporal completeness when using enhanced catalogs and motivates further work into what the spatial covariates of this dataset might offer when combined with flexible point process models such as this one.

## 8 Open Research

The Amatrice-Visso-Norcia catalog produced by Tan et al. (2021b) is accessible at the Zenodo repository <https://doi.org/10.5281/zenodo.4736089> (Tan et al., 2021a). The ETAS simulator used to generate the synthetic data was written for Mizrahi et al. (2021a) and Mizrahi et al. (2021b), and is available at <https://github.com/lmizrahi/etas>, (Mizrahi & Schmid, 2022). Both synthetic and real datasets are found in the reproducibility package along with the models and the experimental design used in this study, found at <https://github.com/ss15859/Neural-Point-Process>.

## Acknowledgments

This project is funded by Compass - Centre for Doctoral Training in Computational Statistics and Data Science (EPSRC Grant Ref EP/S023569/1). Compass is funded by United Kingdom Research and Innovation (UKRI) through the Engineering and Physical Sciences Research Council (EPSRC), <https://www.ukri.org/councils/epsrc>. This project also has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement no. 821115, Real-time earthquake rIsk reduction for a reSilient Europe (RISE), <http://www.rise-eu.org>).

## References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., . . . Zheng, X. (2015). *TensorFlow: Large-scale machine learning on heterogeneous systems*. Retrieved from <http://tensorflow.org/> (Software available from tensorflow.org)
- Bebbington, M., & Harte, D. (2003). The linked stress release model for spatio-temporal seismicity: formulations, procedures and applications. *Geophysical Journal International*, 154(3), 925–946.

- Biloš, M., Sommer, J., Rangapuram, S. S., Januschowski, T., & Günnemann, S. (2021). Neural flows: Efficient alternative to neural odes. *Advances in Neural Information Processing Systems*, 34, 21325–21337.
- Brodsky, E. E. (2011). The spatial density of foreshocks. *Geophysical Research Letters*, 38(10).
- Cattania, C., Werner, M. J., Marzocchi, W., Hainzl, S., Rhoades, D., Gerstenberger, M., ... others (2018). The forecasting skill of physics-based seismicity models during the 2010–2012 canterbury, new zealand, earthquake sequence. *Seismological Research Letters*, 89(4), 1238–1250.
- Chen, R. T., Amos, B., & Nickel, M. (2020). Neural spatio-temporal point processes. *arXiv preprint arXiv:2011.04583*.
- Chiaraluce, L., Di Stefano, R., Tinti, E., Scognamiglio, L., Michele, M., Casarotti, E., ... others (2017). The 2016 central italy seismic sequence: A first look at the mainshocks, aftershocks, and source models. *Seismological Research Letters*, 88(3), 757–771.
- Chilinski, P., & Silva, R. (2020). Neural likelihoods via cumulative distribution functions. In *Conference on uncertainty in artificial intelligence* (pp. 420–429).
- Cleveland, W. S. (1979). Robust locally weighted regression and smoothing scatterplots. *Journal of the American statistical association*, 74(368), 829–836.
- Daley, D. J., & Vere-Jones, D. (2003). Basic properties of the poisson process. *An Introduction to the Theory of Point Processes: Volume I: Elementary Theory and Methods*, 19–40.
- Daley, D. J., Vere-Jones, D., et al. (2003). *An introduction to the theory of point processes: volume i: elementary theory and methods*. Springer.
- Dieterich, J. (1994). A constitutive law for rate of earthquake production and its application to earthquake clustering. *Journal of Geophysical Research: Solid Earth*, 99(B2), 2601–2618.
- Du, N., Dai, H., Trivedi, R., Upadhyay, U., Gomez-Rodriguez, M., & Song, L. (2016). Recurrent marked temporal point processes: Embedding event history to vector. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 1555–1564).
- Ebrahimian, H., & Jalayer, F. (2017). Robust seismicity forecasting based on bayesian parameter estimation for epidemiological spatio-temporal aftershock

- clustering models. *Scientific reports*, 7(1), 1–15.
- Field, E. H., Jordan, T. H., Page, M. T., Milner, K. R., Shaw, B. E., Dawson, T. E., ... others (2017). A synoptic view of the third uniform california earthquake rupture forecast (ucurf3). *Seismological Research Letters*, 88(5), 1259–1267.
- Gulia, L., Rinaldi, A. P., Tormann, T., Vannucci, G., Enescu, B., & Wiemer, S. (2018). The effect of a mainshock on the size distribution of the aftershocks. *Geophysical Research Letters*, 45(24), 13–277.
- Gutenberg, B., & Richter, C. F. (1936). Magnitude and energy of earthquakes. *Science*, 83(2147), 183–185.
- Hainzl, S. (2016a). Apparent triggering function of aftershocks resulting from rate-dependent incompleteness of earthquake catalogs. *Journal of Geophysical Research: Solid Earth*, 121(9), 6499–6509.
- Hainzl, S. (2016b). Rate-dependent incompleteness of earthquake catalogs. *Seismological Research Letters*, 87(2A), 337–344.
- Hainzl, S., Christophersen, A., & Enescu, B. (2008). Impact of earthquake rupture extensions on parameter estimations of point-process models. *Bulletin of the Seismological Society of America*, 98(4), 2066–2072.
- Hawkes, A. G. (1971). Spectra of some self-exciting and mutually exciting point processes. *Biometrika*, 58(1), 83–90.
- Helmstetter, A., Kagan, Y., & Jackson, D. (2007). High-resolution time-independent grid-based forecast for  $m_L = 5$  earthquakes in california. *Seismological Research Letters*, 78(1), 78–86.
- Helmstetter, A., Kagan, Y. Y., & Jackson, D. D. (2006). Comparison of short-term and time-independent earthquake forecast models for southern california. *Bulletin of the Seismological Society of America*, 96(1), 90–106.
- Helmstetter, A., & Sornette, D. (2003). Importance of direct and indirect triggered seismicity in the etas model of seismicity. *Geophysical Research Letters*, 30(11).
- Helmstetter, A., & Werner, M. J. (2014). Adaptive smoothing of seismicity in time, space, and magnitude for time-dependent earthquake forecasts for california. *Bulletin of the Seismological Society of America*, 104(2), 809–822.
- Hochreiter, S. (1998). The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness*



and *Knowledge-Based Systems*, 6(02), 107–116.

Huang, H., Wang, H., & Mak, B. (2019). Recurrent poisson process unit for speech recognition. In *Proceedings of the aaai conference on artificial intelligence* (Vol. 33, pp. 6538–6545).

Kafka, A. L., & Levin, S. Z. (2000). Does the spatial distribution of smaller earthquakes delineate areas where larger earthquakes are likely to occur? *Bulletin of the Seismological Society of America*, 90(3), 724–738.

Kafka, A. L., & Walcott, J. R. (1998). How well does the spatial distribution of smaller earthquakes forecast the locations of larger earthquakes in the north-eastern united states? *Seismological Research Letters*, 69(5), 428–440.

Kagan, Y. Y. (1991). Likelihood analysis of earthquake catalogues. *Geophysical journal international*, 106(1), 135–148.

Kagan, Y. Y. (2004). Short-term properties of earthquake catalogs and models of earthquake source. *Bulletin of the Seismological Society of America*, 94(4), 1207–1228.

Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., & Gerstoft, P. (2019). Machine learning in seismology: Turning data into insights. *Seismological Research Letters*, 90(1), 3–14.

Lawrence, S., Giles, C. L., & Tsoi, A. C. (1997). Lessons in neural network training: Overfitting may be harder than expected. In *Aaai/iaai* (pp. 540–545).

Li, S., Xiao, S., Zhu, S., Du, N., Xie, Y., & Song, L. (2018). Learning temporal point processes via reinforcement learning. *arXiv preprint arXiv:1811.05016*.

Lippiello, E., Marzocchi, W., De Arcangelis, L., & Godano, C. (2012). Spatial organization of foreshocks as a tool to forecast large earthquakes. *Scientific reports*, 2(1), 1–6.

Mancini, S., Segou, M., Werner, M., & Cattania, C. (2019). Improving physics-based aftershock forecasts during the 2016–2017 central italy earthquake cascade. *Journal of Geophysical Research: Solid Earth*, 124(8), 8626–8643.

Mancini, S., Segou, M., Werner, M., Parsons, T., Beroza, G., & Chiaraluce, L. (2022). On the use of high-resolution and deep-learning seismic catalogs for short-term earthquake forecasts: Potential benefits and current limitations.

- Journal of Geophysical Research: Solid Earth*, e2022JB025202.
- Mancini, S., Segou, M., Werner, M. J., & Parsons, T. (2020). The predictive skills of elastic coulomb rate-and-state aftershock forecasts during the 2019 ridge-crest, california, earthquake sequence. *Bulletin of the Seismological Society of America*, 110(4), 1736–1751.
- Marsan, D. (2005). The role of small earthquakes in redistributing crustal elastic stress. *Geophysical Journal International*, 163(1), 141–151.
- Marsan, D., & Lengline, O. (2008). Extending earthquakes’ reach through cascading. *Science*, 319(5866), 1076–1079.
- Marzocchi, W., Lombardi, A. M., & Casarotti, E. (2014). The establishment of an operational earthquake forecasting system in italy. *Seismological Research Letters*, 85(5), 961–969.
- Marzocchi, W., Taroni, M., & Falcone, G. (2017). Earthquake forecasting during the complex amatrice-norcia seismic sequence. *Science advances*, 3(9), e1701239.
- McGuire, J. J., Boettcher, M. S., & Jordan, T. H. (2005). Foreshock sequences and short-term earthquake predictability on east pacific rise transform faults. *Nature*, 434(7032), 457–461.
- Mizrahi, L., Nandan, S., & Wiemer, S. (2021a). The effect of declustering on the size distribution of mainshocks. *Seismological Society of America*, 92(4), 2333–2342.
- Mizrahi, L., Nandan, S., & Wiemer, S. (2021b). Embracing data incompleteness for better earthquake forecasting. *Journal of Geophysical Research: Solid Earth*, 126(12), e2021JB022379.
- Mizrahi, L., & Schmid, N. (2022). *lmizrahi/etas* [Software]. Zenodo. Retrieved from <https://zenodo.org/record/6951562>
- Nandan, S., Ouillon, G., Woessner, J., Sornette, D., & Wiemer, S. (2016). Systematic assessment of the static stress triggering hypothesis using interearthquake time statistics. *Journal of Geophysical Research: Solid Earth*, 121(3), 1890–1909.
- Ogata, Y. (1981). On lewis’ simulation method for point processes. *IEEE transactions on information theory*, 27(1), 23–31.
- Ogata, Y. (1988). Statistical models for earthquake occurrences and residual analysis for point processes. *Journal of the American Statistical association*,

83(401), 9–27.

- Ogata, Y. (1998). Space-time point-process models for earthquake occurrences. *Annals of the Institute of Statistical Mathematics*, 50(2), 379–402.
- Ogata, Y., & Katsura, K. (2014). Comparing foreshock characteristics and foreshock forecasting in observed and simulated earthquake catalogs. *Journal of Geophysical Research: Solid Earth*, 119(11), 8457–8477.
- Omi, T., Ogata, Y., Hirata, Y., & Aihara, K. (2014). Estimating the etas model from an early aftershock sequence. *Geophysical Research Letters*, 41(3), 850–857.
- Omi, T., ueda, n., & Aihara, K. (2019). Fully neural network based model for general temporal point processes. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 32, pp. 2122–2132). Curran Associates, Inc. Retrieved from <https://proceedings.neurips.cc/paper/2019/file/39e4973ba3321b80f37d9b55f63ed8b8-Paper.pdf>
- Rasmussen, J. G. (2018). Lecture notes: Temporal point processes and the conditional intensity function. *arXiv preprint arXiv:1806.00221*.
- Rhoades, D., Liukis, M., Christophersen, A., & Gerstenberger, M. (2016). Retrospective tests of hybrid operational earthquake forecasting models for canterbury. *Geophysical Journal International*, 204(1), 440–456.
- Schorlemmer, D., Wiemer, S., & Wyss, M. (2005). Variations in earthquake-size distribution across different stress regimes. *Nature*, 437(7058), 539–542.
- Seif, S., Mignan, A., Zechar, J. D., Werner, M. J., & Wiemer, S. (2017). Estimating etas: The effects of truncation, missing data, and model assumptions. *Journal of Geophysical Research: Solid Earth*, 122(1), 449–469.
- Shchur, O., Biloš, M., & Günnemann, S. (2019). Intensity-free learning of temporal point processes. *arXiv preprint arXiv:1909.12127*.
- Shchur, O., Türkmen, A. C., Januschowski, T., & Günnemann, S. (2021). Neural temporal point processes: A review. *arXiv preprint arXiv:2104.03528*.
- Tan, Y. J., Waldhauser, F., Ellsworth, W. L., Zhang, M., Zhu, W., Michele, M., . . . Segou, M. (2021a). *Machine-learning-based high-resolution earthquake catalog for the 2016–2017 central italy sequence* [Dataset]. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.4736089>

- Tan, Y. J., Waldhauser, F., Ellsworth, W. L., Zhang, M., Zhu, W., Michele, M., ... Segou, M. (2021b). Machine-learning-based high-resolution earthquake catalog reveals how complex fault structures were activated during the 2016–2017 central italy sequence. *The Seismic Record*, 1(1), 11–19.
- Taroni, M., Marzocchi, W., Schorlemmer, D., Werner, M. J., Wiemer, S., Zechar, J. D., ... Euchner, F. (2018). Prospective csep evaluation of 1-day, 3-month, and 5-yr earthquake forecasts for italy. *Seismological Research Letters*, 89(4), 1251–1261.
- Upadhyay, U., De, A., & Gomez-Rodriguez, M. (2018). Deep reinforcement learning of marked temporal point processes. *arXiv preprint arXiv:1805.09360*.
- Utsu, T. (1955). A relation between the area of after-shock region and the energy of main-shock. *J. Seismol. Soc. Jpn.*, 2, 7, 233–240.
- Utsu, T. (1970). Aftershocks and earthquake statistics (1): Some parameters which characterize an aftershock sequence and their interrelations. *Journal of the Faculty of Science, Hokkaido University. Series 7, Geophysics*, 3(3), 129–195.
- Utsu, T. (1971). Aftershocks and earthquake statistics (2): further investigation of aftershocks and other earthquake sequences based on a new classification of earthquake sequences. *Journal of the Faculty of Science, Hokkaido University. Series 7, Geophysics*, 3(4), 197–266.
- Utsu, T., Ogata, Y., et al. (1995). The centenary of the omori formula for a decay law of aftershock activity. *Journal of Physics of the Earth*, 43(1), 1–33.
- Van Merriënboer, B., Breuleux, O., Bergeron, A., & Lamblin, P. (2018). Automatic differentiation in ml: Where we are and where we should be going. *Advances in neural information processing systems*, 31.
- Wang, Q., Schoenberg, F. P., & Jackson, D. D. (2010). Standard errors of parameter estimates in the etas model. *Bulletin of the Seismological Society of America*, 100(5A), 1989–2001.
- Werner, M. J., Helmstetter, A., Jackson, D. D., & Kagan, Y. Y. (2011). High-resolution long-term and short-term earthquake forecasts for california. *Bulletin of the Seismological Society of America*, 101(4), 1630–1648.
- Wiemer, S., & Wyss, M. (2000). Minimum magnitude of completeness in earthquake catalogs: Examples from alaska, the western united states, and japan. *Bulletin of the Seismological Society of America*, 90(4), 859–869.

- Woessner, J., Hainzl, S., Marzocchi, W., Werner, M., Lombardi, A., Catalli, F., . . .
- Wiemer, S. (2011). A retrospective comparative forecast test on the 1992 landers sequence. *Journal of Geophysical Research: Solid Earth*, 116(B5).
- Woessner, J., & Wiemer, S. (2005). Assessing the quality of earthquake catalogues: Estimating the magnitude of completeness and its uncertainty. *Bulletin of the Seismological Society of America*, 95(2), 684–698.
- Xiao, S., Yan, J., Yang, X., Zha, H., & Chu, S. (2017). Modeling the intensity function of point process via recurrent neural networks. In *Proceedings of the aaai conference on artificial intelligence* (Vol. 31).
- Xiaogu, Z., & Vere-Jones, D. (1994). Further applications of the stochastic stress release model to historical earthquake data. *Tectonophysics*, 229(1-2), 101–121.
- Ying, X. (2019). An overview of overfitting and its solutions. In *Journal of physics: Conference series* (Vol. 1168, p. 022022).
- Zheng, X.-G., & Vere-Jones, D. (1991). Application of stress release models to historical earthquakes from north china. *Pure and Applied Geophysics*, 135(4), 559–576.
- Zhu, W., & Beroza, G. C. (2019). Phasenet: a deep-neural-network-based seismic arrival-time picking method. *Geophysical Journal International*, 216(1), 261–273.
- Zhuang, J., Ogata, Y., & Wang, T. (2017). Data completeness of the kumamoto earthquake sequence in the jma catalog and its influence on the estimation of the etas parameters. *Earth, Planets and Space*, 69(1), 1–12.
- Zhuang, J., Werner, M. J., Hainzl, S., Harte, D., & Zhou, S. (2012). Basic models of seismicity: temporal models. *Community Online Resource for Statistical Seismicity Analysis*.

# Supporting Information for “Forecasting the 2016-2017 Central Apennines Earthquake Sequence with a Neural Point Process”

Samuel Stockman <sup>1</sup>, Daniel J. Lawson <sup>1</sup>, Maximilian J. Werner <sup>2</sup>

<sup>1</sup>School of Mathematics, University of Bristol

<sup>2</sup>School of Earth Sciences, University of Bristol

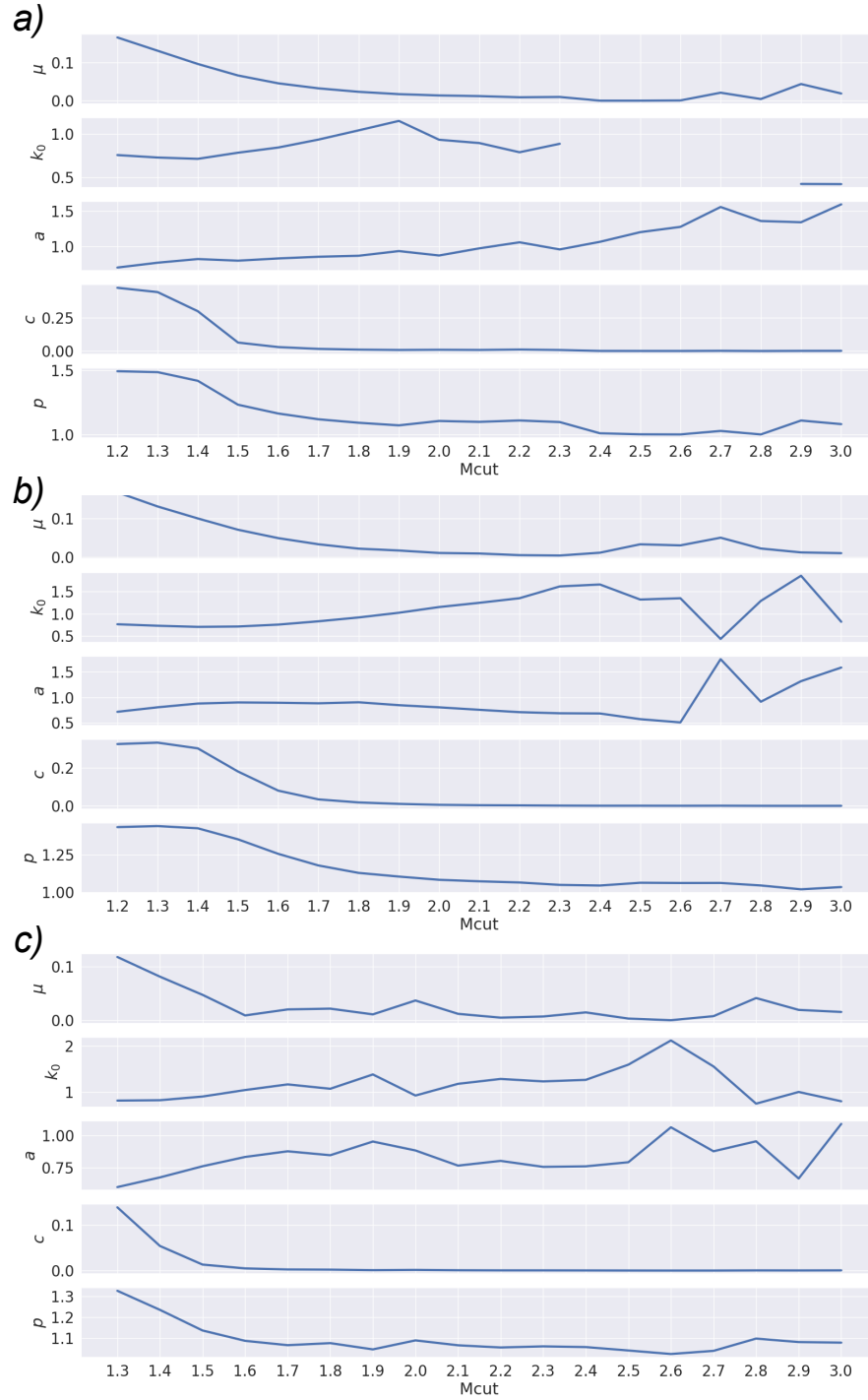
## Contents of this file

1. Figures S1 to S3

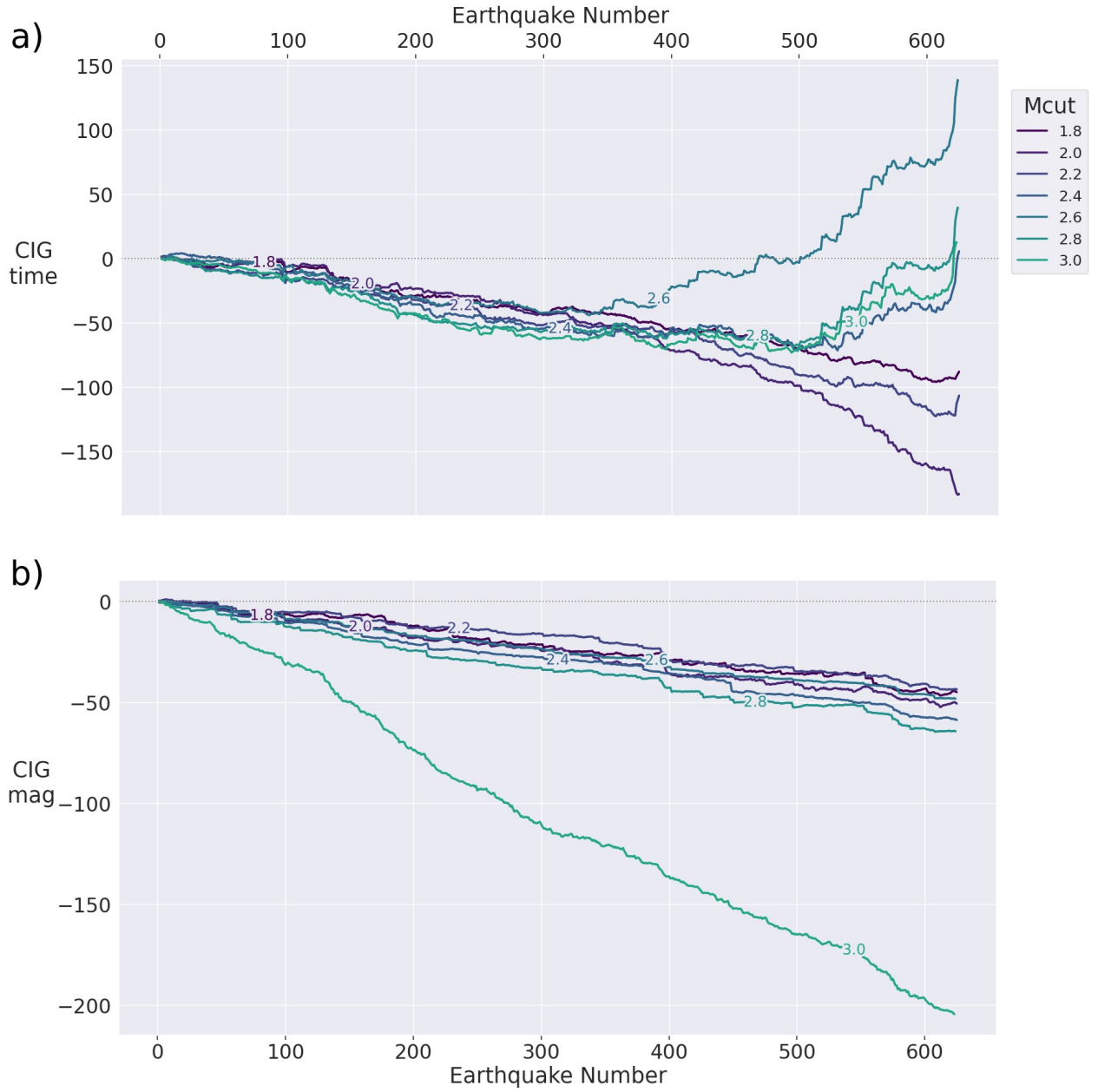
## Introduction

In this supplement, we share the learnt ETAS parameters for the Amatrice-Visso-Norcia catalog for all the values of  $M_{\text{cut}}$  and all training testing partitions (Figure S1). Furthermore we show the cumulative information gain (CIG) of the neural model over ETAS on the complete synthetic catalog for both time and magnitude forecasting (Figure S2), as well as for the incomplete synthetic catalog (Figure S3).

---

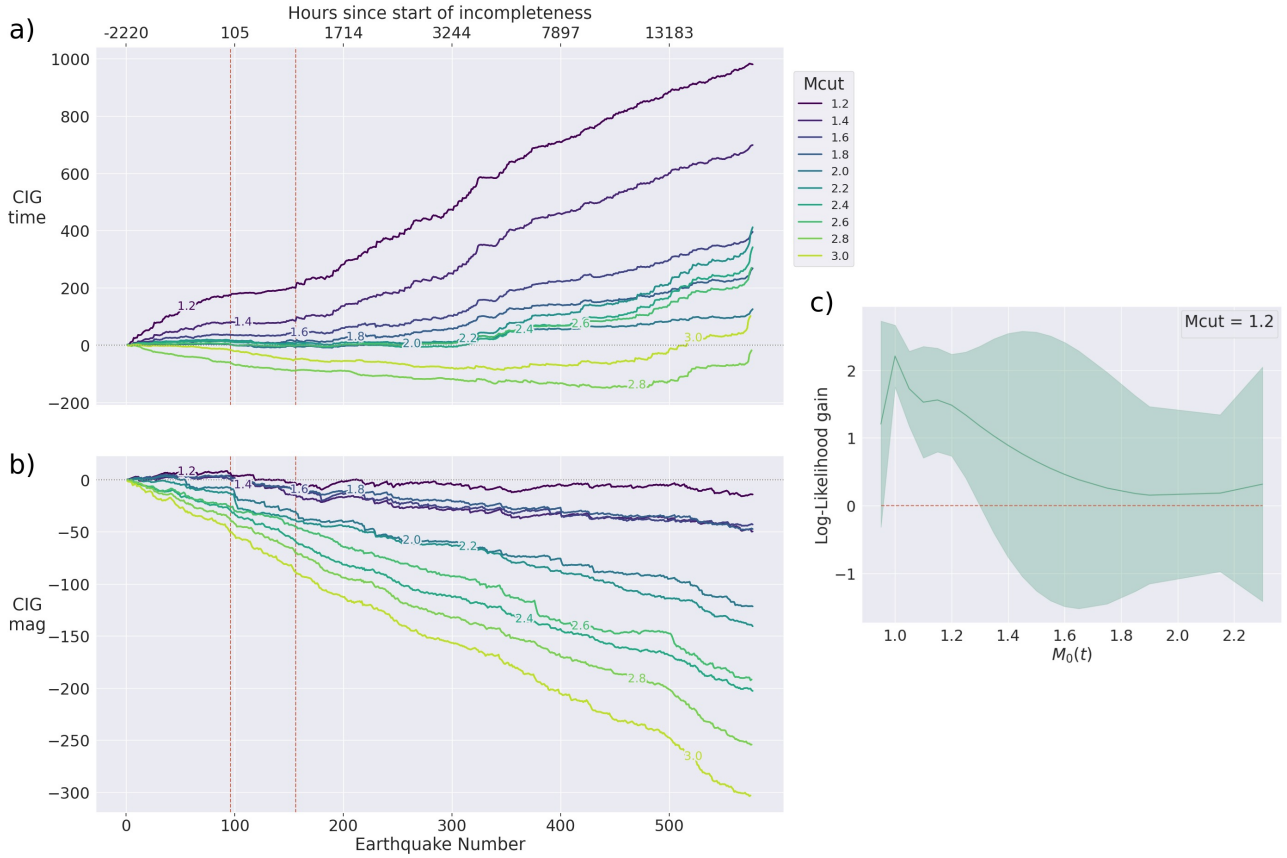


**Figure S1.** Fitted ETAS parameters as a function of  $M_{cut}$ . **a)** training up to the Visso earthquake.  $k_0$  parameters for  $M_{cut}$  2.4-2.8 have been removed from the plot to aid in visualisation. These parameter values are orders of magnitude larger. **b)** training up to the Norcia earthquake. **c)** training up to the Campotosto earthquakes. The unit of time is hours.



**Figure S2.** a) - b) The Cumulative Information Gain (CIG) of the neural model over ETAS for a range of values of  $M_{cut}$ . The models are trained and forecasted on the complete synthetic catalog and the plot depicts the evolution of the CIG from the beginning of the testing period to the end of the catalog. a) displays the CIG for event-time forecasts, b) displays the CIG for magnitude forecasts.





**Figure S3.** a) - b) The Cumulative Information Gain (CIG) of the neural model over ETAS for a range of values of  $M_{cut}$ . The models are trained and forecasted on the incomplete synthetic catalog and the plot depicts the evolution of the CIG from the beginning of the testing period to the end of the catalog. The curve is plotted per event, however, the time since the start of a period of incompleteness is displayed on the top axis. a) displays the CIG for event-time forecasts, b) displays the CIG for magnitude forecasts. c) displays the information gain of the neural model over ETAS as a function of the completeness of the testing catalog - both models are trained with  $M_{cut} = 1.2$ .